

ON PREDICTIVE INFERENCE

by

David Hinkley

University of Minnesota

Technical Report No. 254

January 1976

# ON PREDICTIVE INFERENCE

by

David Hinkley  
University of Minnesota

The likelihood function is the common basis of all parametric inference. However, with the exception of an ad hoc definition by R.A. Fisher, there is no such unifying basis for prediction of future events, given past observations. This article proposes a definition of predictive likelihood which can help to remove some non-uniqueness problems in sampling-theory predictive inference, and which can produce a simple prediction analog of the Bayes result,  $\text{posterior} \propto \text{prior} \times \text{likelihood}$ , in many situations.

AMS 1970 subject classifications.

Key words and phrases. Prediction, confidence regions, exponential families, likelihood, Bayesian methods, pivotal, cross-validation

## 1. Introduction

In 1920 Karl Pearson posed "The Fundamental Problem of Applied Statistics" as follows:

'An "event" has occurred  $p$  times out  $p + q = n$  trials, where we have no a priori knowledge of the frequency of the event in the total population of occurrences. What is the probability of its occurring  $r$  times in a further  $r + s = m$  trials?' Pearson's purpose was to reexamine the general applicability of Bayes's earlier solution, and the resulting controversy, described by Edwards [ 4 ], is of some interest. However, the main question seems to have been largely ignored in the intervening years, while parametric inference has occupied statistical thought. Some authors, for example de Finetti [ 7 ] and Geisser [11], have suggested that more attention be paid to prediction of variables.

Pearson's question could be interpreted in at least two ways. For example, if  $\theta$  is the frequency of the "event" in the total population, so that the unknowable answer to the Fundamental Problem is  $\binom{m}{r} \theta^r (1-\theta)^{m-r}$ , then the question could be answered by giving a point estimate for this real frequency probability. This means interpreting the problem as a decision problem. Alternatively, we might view the problem as an inferential one, as Pearson and Fisher ([16],[17],[18]) did, for which the answer is in terms of relative credibilities for various values of  $r$ . Of course a Bayesian solution does this in terms of a different version of probability, and one might expect that other types of credibility exist for non-Bayesians.

It is to this inferential aspect of prediction that the present paper is directed. The main purpose is to present and discuss a definition

paper is directed. The main purpose is to discuss and discuss a definition  
as to this important aspect of production and the demand  
for non-durable.

On production, and one might expect that other types of expenditure exist  
of course a physical notion does exist in terms of a different notion  
the answer is in terms of relative expenditures for various aspects of a  
production process. In terms of expenditure, the answer is in terms of a  
relative expenditure. This seems intuitively obvious and the problem is to  
show the direction could be answered by having a formal definition for this  
so that the expenditure answer to the expenditure problem is  $\frac{1}{n} \cdot \frac{1}{n-1}$  (1-1)  $\frac{1}{n-1}$   
example, it is the production of the event in the total production  
process's direction could be understood in different ways. For  
of production.

General (1911) have suggested that more attention be paid to production  
expenditure. Some support for example is given by (1921) and  
in the interesting paper, while the expenditure answer has been  
interest. However, the main direction seems to have been largely ignored  
and the following conclusion: (1911) is of some  
importance the general applicability of paper's answer solution.  
It seems in a further  $n + 1 = n$  units. However, the answer is to  
production of consumption. What is the applicability of the economic  
head up a better knowledge of the expenditure of the event in the paper  
the answer has occurred in terms of  $n + 1 = n$  units, where as  
"expenditure" is follows:

In 1920 the answer based on the fundamental problem of production.

1. Introduction

of predictive likelihood, analogous to parametric likelihood, which is Fisherian in spirit but which, strangely, was not suggested by Fisher. The two desirable properties of likelihood, namely resolution of non-unique frequentist confidence regions and ability to yield Bayesian methods, have motivated the definition. Quite by chance, S. Lauritzen and the author proposed much the same likelihood in early 1974, but our work has pursued different paths; Lauritzen's mathematical account has since appeared in [14]. The reader will find the present account deliberately non-mathematical in character; hopefully this will allow easy digestion of the statistical ideas.

Section 2 briefly reviews the two main sampling-theory methods of predictive inference which lead to confidence statements about future random variables, as well as the Bayesian method of defining posterior predictive distributions. The review suggests the need for an analog of parametric likelihood. In Section 3 the analog is given, and shown to be a factor in the Bayes posterior predictive distribution. There is also some discussion of the likelihood as a credibility measure.



## 2. A Brief Critique of Standard Forms of Predictive Inference

### 2.1 Introduction

Suppose that  $Y = (Y_1, \dots, Y_n)$  is a vector random variable with joint probability density function, p.d.f.,  $f_Y(y; \theta)$ ,  $\theta \in \Omega$ , and that the random variable  $Z = (Z_1, \dots, Z_m)$  independently has p.d.f.  $g_Z(z; \theta)$ . An observation  $y$  on  $Y$  is available, and we wish to predict  $Z$  in the absence of knowledge about  $\theta$ . The independence of  $Y$  and  $Z$  is assumed only for simplicity, and the general theory outlined below extends easily.

Prediction of  $Z$  could mean many things. It could mean estimation of the p.d.f.  $g_Z$ , or its probability integral. Or, prediction could mean construction of a series of confidence regions. In the absence of a utility specification, this latter interpretation is the usual one, which we discuss here. First, in Section 2.2, the two main sampling-theory types of confidence region (critical region and pivotal) are described, and their non-uniqueness emphasized. Section 2.3 defines the Bayesian posterior predictive distribution and suggests the desirability of a prediction analog of parametric likelihood.

### 2.2 Sampling-theory Predictive Inference

A confidence region for  $Z$  is a set  $P_\alpha(Y)$  in the sample space of  $Z$ , determined by the observable  $Y$  and satisfying

$$\text{pr}\{Z \in P_\alpha(Y); \theta\} = 1 - \alpha \quad (2.1)$$

independently of  $\theta$ ; several regions, for different values of  $\alpha$ , might be given simultaneously. There are two main methods of constructing  $P_\alpha(Y)$ , one based on test critical regions, the other based on pivotal quantities.

Construction of  $P_\alpha(Y)$  via test critical regions proceeds as follows: Suppose that  $Z$  has p.d.f.  $g_Z(z; \theta^*)$ , and consider the hypothesis  $H_0: \theta = \theta^*$

$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n f\left(\frac{k}{n}\right) = \int_0^1 f(x) dx$

consideration of  $E^U(X)$  are also sufficient to prove the following:  
the proof of the existence of the above mentioned structures.  
which are the only structures of consideration  $E^U(X)$ .  
the proof of the existence of the above mentioned structures of the above mentioned structures.

$$\text{def } \text{is\_prime}(x): \text{ return } x \in \{2, 3, 5, 7, 11, 13, 17, 19, 23, 29, 31, 37, 41, 43, 47, 53, 59, 61, 67, 71, 73, 79, 83, 89, 97\} \quad (3.7)$$

SECRETORUL SA FUE ORDONATO SA FIE ESTIMAT.

CONFIDENTIAL

SECRET - SECURITY INFORMATION

REGISTRATION STATE OF DELAWARE

[illegible]

and upon our-  
self.

phases of conjugated action (couples, action and synthesis) are described:

15. Средства связи: телефон ИП 800-200-10-10 для договора субсидиарной

[illegible]

consideration of a series of contingencies before the release of a

of the B.G.P. 300 or the biographical type. The biographical type is

ΣΥΝΟΨΗ ΤΩΝ ΣΥΝΤΑΚΤΙΚΩΝ ΚΑΙ ΣΥΝΤΑΞΕΩΝ ΤΗΣ ΕΛΛΗΝΙΚΗΣ ΓΡΑΜΜΑΤΙΚΗΣ

14-00000

01 июль 1986 года - при переселении из п. Ан-1 в поселок Ошанька.

QD86L4870N A ON Y FO SASTSNG'S' SDC AS ATST JO BUCONOR P IN DUS BOBONOS

REASONABLE ASSETS IN THE HANDS OF THE PERSONS WHOSE NAMES ARE LISTED IN THE ATTACHED LIST OF NAMES.

SECRET

2.  $\bar{A} = (A^1, \dots, A^N)$  is a fixed random vector with

# DECLASSIFICATION

7. BIRTH RECORDING OF CRIMINALS: POLICE OF MICHIGAN THROUGHOUT



with some alternative hypothesis  $H_A$ . If for the particular  $H_A$  a size -  $\alpha$  similar test critical region can be found, then by definition

$$\text{pr}\{(Y, Z) \in R_\alpha; H_0\} = \alpha. \quad (2.2)$$

But  $H_0$  is true, and so if  $P_\alpha(y)$  is the projection of the complement  $\bar{R}_\alpha$  onto  $Y = y$ , then  $P_\alpha(Y)$  satisfies the requirement (2.1). Thus the prediction confidence region contains all values of  $z$  which are not in conflict with the model and observation  $y$ , as judged by a size -  $\alpha$  test.

It is clear that many such  $R_\alpha$  may exist, and their existence is not dependent on the nature of  $H_A$ . However, the standard procedure is to determine a "good" critical region for the particular alternative  $H_A$  introduced, in the sense of being most powerful (locally or uniformly). For many problems this results in precise confidence statements about  $Z$ , but the general concept of power seems somewhat remote from prediction.

Two important points to note about this approach are:

- (i) If the minimal sufficient statistic under  $H_0$  is not a reduction of the minimal sufficient statistics for the two families, then the Neyman structure property of similar tests is lost, and there may be no similar critical regions  $R_\alpha$ . A simple example is that where the observables  $Y_1, \dots, Y_n$  are independent counts from a Poisson process of rate  $\theta$ , and  $Z$  is the time from "now" to the occurrence of the next event.
- (ii) The method is not uniquely defined in the sense that different alternatives  $H_A$  are possible. Indeed if  $\dim(\theta) > \dim(Z)$ , then there must be a restriction on  $H_A$  in order to obtain non-trivial critical regions, and there is no unique restriction. For example, if  $Y_1, \dots, Y_n$ , are i.i.d.

under restriction. For example, let  $H^1, \dots, H^k$  are i.i.d.

in order to obtain non-trivial critical regions, and there is no  $\dim(S) > \dim(V)$ , then there must be a restriction on  $H^i$

different alternatives  $H^i$  may be chosen. Indeed let  $\dim(S)$

(iii) The method is not directly defined in the sense that the time from "now" to the occurrence of the next event.

Independent counts from a Poisson process of rate  $\lambda$  and  $\mu$  is simple example is that there are observations  $X^1, \dots, X^k$  are

is lost, and there may be no similar critical regions  $H^1, \dots, H^k$  families, then the Nelson structure hypothesis of similar cases

reduction of the minimal sufficient statistic for the two

(1) is the minimal sufficient statistic region  $H^0$  is not a

two important points to note about this approach are:

But the general concept of power seems somewhat remote from prediction.

For many problems this results in precise, complete or statements about  $N^i$  introduced, in the sense of being most powerful (locally or uniformly).

determining a "good" critical region for the restricted alternatives  $H^i$  dependent on the nature of  $H^i$ . However, the standard procedure is to

is to check first many such  $H^i$  and select, and check evidence is not conflict with the model and observation  $X^i$  is judged by a given test.

prediction confidence region contains all values of  $X$  which are not in

$H^0$  or  $H^1 = X^1$  then  $H^1(X)$  satisfies the requirement (C.1). Thus the

but  $H^0$  is true, and so  $H^1(X)$  is the projection of the complement

$$\text{But } (X^1, X^2) \in H^1 : H^0 = \emptyset. \quad (S.5)$$

similar test critical region can be found, then the restriction

with some alternative hypothesis  $H^i$ . In fact the predictor  $H^i$ 's size--

$N(\mu, \tau)$  and  $Z_1$  is independently  $N(\mu^*, \tau^*)$ , then we cannot allow both  $\mu \neq \mu^*$  and  $\tau \neq \tau^*$ ; the form of the prediction confidence region will depend on  $H_A$  if a most powerful critical region is used.

The above points are made to imply a degree of arbitrariness in the critical region construct. Similar criticisms can be made against critical region construction of confidence regions for  $\theta$ , but in that case non-uniqueness is removable by use of the likelihood function. That is to say, among the systems of confidence regions  $C_\alpha(Y)$  for  $\theta$  which exist with the required property

$$\text{pr}\{\theta \in C_\alpha(Y); \theta\} = 1 - \alpha ,$$

we take the uniquely-defined likelihood-based confidence region  $C_\alpha^*(y)$  satisfying

$$\inf_{\theta' \in C_\alpha^*(y)} \text{lik}(\theta' | y) \geq \sup_{\theta'' \notin C_\alpha^*(y)} \text{lik}(\theta'' | y) . \quad (2.3)$$

It seems a reasonable requirement that no parameter point outside the confidence region have higher likelihood than any point inside the region, and a corresponding likelihood base for our predictand  $Z$  is much to be desired.

The other main method for obtaining prediction confidence regions is the use of pivotal quantities. This has strong connections with fiducial and structural inference, but we shall not pursue them here. Briefly, a pivotal quantity in the prediction context is a function  $r(Y, Z)$  with distribution independent of  $\theta$ . Thus for any  $\alpha$ , we can in principle determine a region  $T_\alpha$  such that

$$\text{pr}\{r(Y, Z) \in T_\alpha; \theta\} = 1 - \alpha \quad (2.4)$$

$$\text{Def } \lambda(A, \alpha) = \lambda(A, \alpha) = \gamma - \alpha$$

(3.2)

Definition 3.2. Let  $\lambda(A, \alpha)$  be a function

$\lambda(A, \alpha)$  after differentiation with respect to  $\alpha$ . Then for any  $\alpha$  we can find a function  $\lambda(A, \alpha)$  in the neighborhood of  $\alpha$  as a function of  $\alpha$  and  $\lambda(A, \alpha)$  is a function of  $\alpha$  and  $\lambda(A, \alpha)$  is a function of  $\alpha$ . Then we can find a function of  $\alpha$  and  $\lambda(A, \alpha)$  is a function of  $\alpha$ . Then we can find a function of  $\alpha$  and  $\lambda(A, \alpha)$  is a function of  $\alpha$ .

The other way round for obtaining a function of  $\alpha$  and  $\lambda(A, \alpha)$  is a function of  $\alpha$  and  $\lambda(A, \alpha)$  is a function of  $\alpha$ .

Definition 3.3. Let  $\lambda(A, \alpha)$  be a function of  $\alpha$  and  $\lambda(A, \alpha)$  is a function of  $\alpha$ . Then we can find a function of  $\alpha$  and  $\lambda(A, \alpha)$  is a function of  $\alpha$ . Then we can find a function of  $\alpha$  and  $\lambda(A, \alpha)$  is a function of  $\alpha$ .

$$\lambda(A, \alpha) = \lambda(A, \alpha) = \gamma - \alpha$$

(3.3)

Definition 3.4

Let  $\lambda(A, \alpha)$  be a function of  $\alpha$  and  $\lambda(A, \alpha)$  is a function of  $\alpha$ .

$$\text{Def } \lambda(A, \alpha) = \lambda(A, \alpha) = \gamma - \alpha$$

Definition 3.5. Let  $\lambda(A, \alpha)$  be a function of  $\alpha$  and  $\lambda(A, \alpha)$  is a function of  $\alpha$ .

Definition 3.6. Let  $\lambda(A, \alpha)$  be a function of  $\alpha$  and  $\lambda(A, \alpha)$  is a function of  $\alpha$ . Then we can find a function of  $\alpha$  and  $\lambda(A, \alpha)$  is a function of  $\alpha$ . Then we can find a function of  $\alpha$  and  $\lambda(A, \alpha)$  is a function of  $\alpha$ .

Definition 3.7. Let  $\lambda(A, \alpha)$  be a function of  $\alpha$  and  $\lambda(A, \alpha)$  is a function of  $\alpha$ . Then we can find a function of  $\alpha$  and  $\lambda(A, \alpha)$  is a function of  $\alpha$ . Then we can find a function of  $\alpha$  and  $\lambda(A, \alpha)$  is a function of  $\alpha$ .

for all  $\theta$ ; cf. (2.2). Then a pivotal prediction region is the projection of  $T_\alpha$  onto  $Y = y$ ,

$$P_\alpha(y) = \{z | r(y,z) \in T_\alpha\} . \quad (2.5)$$

This approach requires construction of an appropriate pivotal quantity  $r(Y, Z)$ , which can be accomplished by an extension of G.A. Barnard's work on pivotal parametric inference; see [1].

We retain the assumption of independence of  $Y$  and  $Z$  for simplicity. Then suppose that the one-one transformation  $\gamma$  takes  $Y$  into

$$C = \begin{pmatrix} C_1 \\ C_2 \end{pmatrix} = \begin{pmatrix} c_1(Y, \theta) \\ c_2(Y) \end{pmatrix} ,$$

where  $C_1$  and  $C_2$  are pivots having distributions independent of  $\theta$ ;  $C_1$  is termed proper, and  $C_2$  is termed ancillary. Suppose that a corresponding transformation  $\delta$  takes  $Z$  into

$$D = \begin{pmatrix} D_1 \\ D_2 \end{pmatrix} = \begin{pmatrix} d_1(Z, \theta) \\ d_2(Z) \end{pmatrix} ,$$

where  $D_1$  and  $D_2$  are respectively proper and ancillary pivotals. Since  $D_2$  has known distribution, the prediction problem focusses on  $D_1$ , to which the information about  $\theta$  in the conditional distribution of  $C_1$  given  $C_2 = c_2$  relates. Now we make the definition of predictive pivot: if a one-one transformation of  $(C_1, D_1)$  exists such that

$$\begin{pmatrix} C_1 \\ D_1 \end{pmatrix} \rightarrow \begin{pmatrix} Q \\ R \end{pmatrix} = \begin{pmatrix} q(Y, Z, \theta) \\ r(Y, Z) \end{pmatrix}$$

then  $R$  is a predictive pivot. With the transformation of  $(Y, Z)$  into  $(C_2, D_2, Q, R)$  is associated the re-expression of their joint distribution,

which takes the form

$(c^{\bar{J}}, p^{\bar{J}})$  is associated the re-expression of such logic transformation:  
 such  $\lambda$  is a bijective block. When the transformation of  $(\lambda^{\bar{J}}, \lambda)$  into

$$\begin{pmatrix} c^{\bar{J}} \\ p^{\bar{J}} \end{pmatrix} \rightarrow \begin{pmatrix} \lambda \\ \lambda \end{pmatrix} = \begin{pmatrix} \lambda(\lambda^{\bar{J}}, \lambda) \\ \lambda(\lambda^{\bar{J}}, \lambda) \end{pmatrix}$$

is a one-one transformation of  $(c^{\bar{J}}, p^{\bar{J}})$  into such that  
 block  $c^{\bar{J}} = c^{\bar{J}}$  remains. Now we have the definition of bijective block:  
 which the transformation from such the complete transformation of  $c^{\bar{J}}$   
 $p^{\bar{J}}$  has known descriptions: the bijective block consists on  $p^{\bar{J}}$ , so  
 blocks  $p^{\bar{J}}$  and  $p^{\bar{J}}$  are bijectively blocks and mutually bijective. Thus,

$$p = \begin{pmatrix} c^{\bar{J}} \\ p^{\bar{J}} \end{pmatrix} = \begin{pmatrix} c^{\bar{J}}(\lambda) \\ c^{\bar{J}}(\lambda^{\bar{J}}, \lambda) \end{pmatrix}.$$

transformation  $\lambda$  takes  $\lambda$  into

is called block, and  $c^{\bar{J}}$  is called subblock. Subblock and a subblock  
 blocks  $c^{\bar{J}}$  and  $c^{\bar{J}}$  are blocks relative transformation transformation of  $c^{\bar{J}}$

$$c = \begin{pmatrix} c^{\bar{J}} \\ c^{\bar{J}} \end{pmatrix} = \begin{pmatrix} c^{\bar{J}}(\lambda) \\ c^{\bar{J}}(\lambda^{\bar{J}}, \lambda) \end{pmatrix}.$$

When subblock and the one-one transformation  $\lambda$  takes  $\lambda$  into

We denote the transformation of transformation of  $\lambda$  and  $\lambda$  for transformation.  
 block on block bijective transformation: see relation (1.12).  
 $\lambda(\lambda^{\bar{J}}, \lambda)$ , which can be associated of an expression of  $\lambda$  and  $\lambda$  relative  
 this expression relative transformation of an expression bijective transformation

$$\lambda^{\bar{J}}(\lambda) = \{\lambda | \lambda(\lambda) \lambda\} = \lambda^{\bar{J}}. \quad (3.2)$$

on  $\lambda^{\bar{J}}$  onto  $\lambda = \lambda^{\bar{J}}$

Now we have: on (3.2): block a bijective transformation relative is the bijective

which takes the form

$$f_Y(y;\theta)g_Z(z;\theta) = \varphi(c_2) \psi(d_2) \xi(r|c_2, d_2) \eta(q|c_2, d_2, r; \theta).$$

The prediction problem is then resolved by discarding  $q(Y, Z, \theta)$ , either on the principle that  $\theta$  and  $Z$  are confounded, or by appeal to lack of invariance of  $Q$  under isomorphic transformations of  $Y$  and  $Z$  which leave  $R$  invariant. With  $Q$  discarded, we are left with the following predictive information:

$c_2(Z)$  has p.d.f.  $\varphi(c_2)$ , and conditional on  $(C_2, D_2) = (c_2, d_2)$  the pivot  $r(Y, Z)$  has p.d.f.  $\xi(r|c_2, d_2)$ , the observed value of  $Y$  being  $y$ .

Note that this leaves ambiguous the particular choice of region  $T_\alpha$ , and hence  $P_\alpha(y)$ , in (2.4). For example, in a location parameter problem where the predictive pivot  $R$  may be chosen as  $\bar{Y} - \bar{Z} = n^{-1} \sum Y_j - m^{-1} \sum Z_j$ ,  $P_\alpha(y)$  is any invariant region of size  $1 - \alpha$  with respect to the conditional distribution of  $R$ , shifted by amount  $\bar{y}$ . The ambiguity would be removed by defining  $\xi(r|c_2, d_2) \varphi(c_2)$  as the likelihood function of  $z$  given  $y$ , to be used to uniquely define  $P_\alpha^*(y)$  in a manner analogous to  $C_\alpha^*(y)$  in (2.3). This definition of likelihood is unsatisfactory on two counts: Firstly, since the pivotal approach fails for most discrete distributions, we would be left with no definition of likelihood. Second, this definition of likelihood would not generally permit combination with prior information about  $Z$  to obtain Bayes posterior predictive distributions; in fact explicit prior information needs to be represented by the pivotal quantity  $@$ , which together with  $C_1$  and  $D_1$  defines the conditional distribution of  $Z$ ; see [1].

1

THOSE WITH THIS POLICE CERTIFICATE ARE DESIGNATED OFFICER OF RECORD

$$= (c^{\beta} \cdot c^{\gamma}) \text{ при } b \cdot c \cdot e = (1^{\beta} \cdot e) \text{ при } b \cdot c \cdot e = (1 \cdot c^{\beta} \cdot e^{\gamma}) \text{ при } \\ c^{\beta} \cdot e \text{ при } b \cdot c \cdot e = (c^{\beta}) \text{ при } b \cdot c \cdot e \text{ или } (c^{\beta} \cdot e^{\gamma}) =$$

Y THIRTYFIVE. FOUR 9 QUESSENG. AS HIS FOUR FOUR TWO TWO

$$T^{\dagger}(z_1) \cdot T^{\dagger}(z_2) = (c^{\dagger}_1 - c^{\dagger}_2) \cdot (c^{\dagger}_2 - c^{\dagger}_1) = (c^{\dagger}_1 \cdot c^{\dagger}_2 - c^{\dagger}_1 \cdot c^{\dagger}_1 - c^{\dagger}_2 \cdot c^{\dagger}_2 + c^{\dagger}_2 \cdot c^{\dagger}_1) = (c^{\dagger}_1 \cdot c^{\dagger}_2 - c^{\dagger}_2 \cdot c^{\dagger}_1) = 0.$$

- 4 -



To sum up: The two most useful approaches to construction of similar prediction confidence regions do not in fact uniquely define such regions, nor do they give solutions to all problems. Partial resolution of these difficulties requires definition of a suitably general predictive likelihood.

### 2.3 Bayesian Predictive Inference

The preceding critique of sampling-theory predictive inference would usually fall on deaf Bayesian ears. However, a prediction analog of parametric likelihood might be of some interest within the Bayesian framework. Recall that for the general case, if  $p_{\Theta}(\theta)$  is the prior density for  $\Theta$ , then the Bayes posterior predictive distribution will have density

$$\begin{aligned} f_{Z|Y}(z|y) &= \int f_{\Theta|Y}(\theta|y) f_{Z|Y,\theta}(z|y,\theta) d\theta \\ &= \frac{\int f_{Y|\Theta}(y|\theta) p_{\Theta}(\theta) f_{Z|Y,\theta}(z|y,\theta) d\theta}{\int f_{Y|\Theta}(y|\theta) p_{\Theta}(\theta) d\theta} . \end{aligned} \quad (2.5)$$

This appears to be somewhat complicated when compared to the "loose" form of the Bayes posterior for  $\Theta$ , i.e.

$$f_{\Theta|Y}(\theta|y) \propto \text{lik}_Y(\theta|y) p_{\Theta}(\theta) .$$

It would, then, be of some interest to obtain a prediction analog of the likelihood  $\text{lik}_Y(\theta|y)$  for which (2.5) is expressible as

$$f_{Z|Y}(z|y) \propto \text{likelihood} \times p_Z(z) , \quad (2.6)$$

where  $P_Z(z)$  is the marginal prior density of  $Z$ . Indeed in many ways (2.6) would seem to be a preferred form of (2.5), since presumably subjective probability assessments concerning  $z$  would often be more tractable than corresponding assessments concerning  $\theta$ .

principles from corresponding assessments concerning a.

subjective probability assessments concerning a world often be more

(5.6) would seem to be a Bayesian form of (3.2), since presumably

where  $P^A(x)$  is the marginal prior density of  $x$ . Placed in this case

$$P^A(x|A) = P^A(x|A)P^A(x) \quad (5.7)$$

where  $P^A(x|A)$  is the marginal prior density of  $x$  for which (3.2) is adequate to be

where  $P^A(x)$  is the marginal prior density of  $x$ . Placed in this case

$$P^A(x|A) = P^A(x|A)P^A(x) \quad (5.8)$$

form of the Bayes posterior for  $x$ , i.e.

will appear to be somewhat complicated when compared to the "proper"

$$= \frac{P^A(x|A)P^A(x|A)}{P^A(x|A)P^A(x|A)} \quad (5.9)$$

$$P^A(x|A) = P^A(x|A)P^A(x|A)$$

will have constant

the prior density for  $x$  from the Bayes decision-theoretic perspective (interpretation  
 Bayesian framework). Locally this for the "general case" of  $P^A(x)$  is  
 instead of Bayesian likelihood which is of some interest which the  
 world usually rely on best Bayesian case. However, a prediction

The preceding critique of sampling-theory predictive inference

5.3 Predictive Bayesian Inference

General predictive hypothesis. Such a likelihood should also be a  
 resolution of these difficulties requires definition of a subjective  
 such regions, not as such given solutions to all problems. Subjective  
 various prediction confidence regions do not in fact uniformly require

to sum up: The two most notable approaches to construction of

### 3. Predictive Likelihoods

#### 3.1 Introduction

The previous section suggests a need for the prediction analog of parametric likelihood. One such analog was introduced, rather naturally, by R.A. Fisher [10], although his use was rather ad hoc. Section 3.2 describes Fisher's definition and shows it to be defective in terms of reasonable consistency requirements. A new definition of predictive likelihood is proposed in Section 3.3, which details various important properties of the likelihood for exponential families, including the satisfaction of the Bayes factorization (2.6). Section 3.4 briefly discusses a conditional version of predictive likelihood that is operational for location-scale families, for example.

#### 3.2 Fisher's Predictive Likelihood

In his book Statistical Methods and Scientific Inference Fisher discussed the "Fundamental Problem" described in Section 1. He proposed a likelihood for the future sample based on the likelihood ratio measure of support for the true hypothesis of constant probability that is obtained from comparison of the observed sample and the future sample. Specifically, if the event has occurred  $y$  times out of  $n$  trials, where the chance of occurrence is  $\theta$  on each trial, and if  $z$  is the number of occurrences in  $m$  further trials, then the predictive likelihood is the likelihood ratio statistic for testing that  $\theta$  has the same value in both samples, namely

$$\begin{aligned} \text{FLIK}(z|y) &= \frac{\sup_{\theta} \binom{m+n}{y+z} \theta^{y+z} (1-\theta)^{m+n-(y+z)}}{\sup_{\varphi} \binom{n}{y} \varphi^y (1-\varphi)^{n-y} \sup_{\varphi} \binom{m}{z} \varphi^z (1-\varphi)^{m-z}} \\ &= \frac{(y+z)^{y+z} (m+n-y-z)^{m+n-y-z} \binom{n}{y} \binom{m}{z}}{y^y (n-y)^{n-y} z^z (m-z)^{m-z}} \end{aligned} \quad (3.1)$$



The idea is clearly similar to that underlying the critical region construction of prediction confidence regions (Section 2.2). The extension to general families of distributions is obvious.

From the form of the binomial result (3.1) it is evident that Fisher's definition will not produce the type of factorization (2.6) for Bayes posterior distributions. (Of incidental interest in this regard is the fact that Fisher [9, p.393] implicitly claimed that such a factorization was impossible.) However, the definition of FLIK leads to unsatisfactory results at a more primitive level, as we shall now show.

One fairly reasonable requirement of predictive likelihood is that infinite amounts of sample data should lead to the actual probability density of the future variable  $Z$ . Of course this is not a logical requirement. In the binomial example, it is quite easy to see that as  $n \rightarrow \infty$ , with  $y/n \rightarrow \theta$ ,

$$\text{FLIK}(z|y) \rightarrow \frac{\binom{m}{z} \theta^z (1-\theta)^{m-z}}{\sup_z \binom{m}{z} \theta^z (1-\theta)^{m-z}},$$

which is the probability density at  $z$  divided by its supremum over the parameter space. This is easy to prove for general regular distributions of  $Y$  and  $Z$ , and so leads to a straightforward modification of Fisher's definition

$$\text{FLIK}^*(z|y) = \frac{\sup_{\theta} \text{lik}_{Y,Z}(\theta; y, z)}{\sup_{\theta} \text{lik}_Y(\theta; y)}. \quad (3.2)$$

This does not resolve matters completely, however, since a second, even more appealing, requirement is that if in the binomial case  $m \rightarrow \infty$  with  $z/m \rightarrow \psi$ , then the predictive likelihood should converge to the ordinary likelihood of  $\psi$  given  $Y = y$ . It is a simple exercise to show that neither (3.1) nor the modification (3.2) possess this property.

Further discussion of Fisher's predictive likelihood may be found in an essay on the "Fundamental Problem" by Edwards [4]. The modification (3.2) has been studied by Faulkenberry and Lejeune (unpublished report, Oregon State University).

### 3.3 A New Predictive Likelihood

The basic idea behind Fisher's proposal was that predictive likelihood should reflect the degree to which the values  $y$  and  $z$  of the observed and unobserved variables support the true hypothesis of common distribution. This idea can be applied in a different way in certain simplifying circumstances, namely when the distributions involved belong to exponential families.

Suppose that  $Y_1, \dots, Y_N$ ,  $N = n + n_F$ , are i.i.d. random variables each with p.d.f.  $f_Y(y; \theta)$ ,  $\theta \in \Omega$ , such that  $Y = (Y_1, \dots, Y_n)$  is to be observed and  $Z = (Y_{n_P+1}, \dots, Y_N)$  is to be predicted. Further suppose that  $s(\cdot)$  is the minimal sufficient reduction function, and denote  $s(Y) = s_P$ ,  $s(Z) = s_F$ ,  $s(Y, Z) = s_{PF}$ . Then we propose to measure the likelihood of  $z$  by the plausibility of the observed value  $y$  given  $s_{PF}$ . This notion is exactly that used in defining parametric likelihood, where  $\theta$  is substituted for  $s_{PF}$ . Specifically, the likelihood is to be proportional to the conditional density of  $Y = y$  given  $S_{PF} = s(y, z)$ . A "surprising" or "unlikely" value of  $z$  will be one which would a posteriori make the sample value  $y$  appear improbable. Since the conditional density of  $Y$  given  $S_{PF}$  factorizes into  $f_{Y|S_P}(y|s(y))f_{S_P|S_{PF}}(s(y)|s(y, z))$ , only the second factor is of interest, so that we have

Definition 1 For i.i.d. random variables  $Y_1, \dots, Y_N$ ,  $N = n_P + n_F$ , such that  $S_P = s(Y_1, \dots, Y_{n_P})$  and  $S_{PF} = s(Y_1, \dots, Y_N)$  are minimal sufficient reductions, the predictive likelihood of  $S_F = s(Z) = s(Y_{n_P+1}, \dots, Y_N)$  is

$$\text{PLIK}(s_F|s_P) = f_{S_P|S_{PF}}(s(y)|s(y, z)) \quad (3.3)$$



Of course this likelihood is independent of  $\theta$  since  $S_{PF}$  is sufficient. Also, given  $S_F$ , the probability distribution of  $Z$  is known, independent of  $\theta$ , so that only the value of  $s_F$  is of concern to us.

Certain interesting features of the definition are clearly illustrated in the case of the "Fundamental Problem".

Example 3.1 Let  $Y_1, \dots, Y_N$  be i.i.d. Bernoulli variables with

$$\text{PR}(Y_j = 1) = \theta = \frac{n}{n_F} \text{PR}(Y_j = 0) .$$

Then if  $n_P$  and  $n_F$  are fixed,  $S_P = \sum_{j=1}^{n_P} Y_j$ ,  $S_F = \sum_{j=1}^{n_F} Y_{n_P+k}$  and  $S_F = S_P + S_F$ . Simple calculation of the conditional probability in Definition 1 gives

$$\text{PLIK}(s_F | s_P) = \frac{\binom{n_P}{s_P} \binom{n_F}{s_F}}{\binom{n_P + n_F}{s_P + s_F}}, \quad (3.4)$$

which is the hypergeometric likelihood for sampling without replacement  $n_P$  times from an urn containing  $s_P + s_F$  "ones" and  $N - (s_P + s_F)$  "zeros." The special case  $n_F = 1$  gives likelihood odds

$$\frac{\text{PLIK}(1 | s_P)}{\text{PLIK}(0 | s_P)} = \frac{s_P + 1}{n_P - s_P + 1} ,$$

formally the same as Laplace's Law of Succession.

The consistency properties described in Section 3.2 are clearly satisfied by (3.4), i.e.

$$(i) \quad \text{if } n_P \rightarrow \infty \text{ such that } n_P^{-1} s_P \rightarrow \theta, \quad n_F \text{ fixed,}$$

$$\lim_{n_P \rightarrow \infty} \text{PLIK}(s_F | s_P) = \binom{n_F}{s_F} \theta^{s_F} (1 - \theta)^{n_F - s_F} = \text{PR}(S_F = s_F | \theta); \quad (3.5)$$

convergence here is w.p.1, at the rate  $n_P^{-1/2}$ .

$$(ii) \quad \text{if } n_F \rightarrow \infty \text{ such that } n_F^{-1} s_F \rightarrow \psi, \quad n_P \text{ fixed,}$$

$$\lim_{n_F \rightarrow \infty} \text{PLIK}(s_F | s_P) = \binom{n_P}{s_P} \psi^{s_P} (1 - \psi)^{n_P - s_P} = \text{lik}_{s_P}(\psi | s_P, n_P). \quad (3.6)$$

□ □ □



$$\text{ITU-BRUK}(e^a | e^b) = \begin{pmatrix} e^b \\ e^b \\ e^b \end{pmatrix} \quad \text{and} \quad \text{ITU-BRUK}(e^a | e^b) = \text{ITU-BRUK}(e^b | e^a) \quad (3.2)$$

CONACILHOS NOME DE A.B.T. e: 2000 47.

$$E_{\text{eff}}(e^{\pm}) = \begin{pmatrix} e^{\pm} \\ e^{\pm} \end{pmatrix} \quad (3.2)$$

(T)  $\mu_{\text{H}}^{\text{H}} = \frac{1}{2}(\mu_{\text{H}}^{\text{H}} + \mu_{\text{H}}^{\text{H}})$

до содержания, подлежащего рассмотрению в пункте 3.5 его отчета.

$$\frac{\text{EFTG}(a^b)}{\text{EFTG}(1/a^b)} = \frac{a^b}{1/a^b} = a^b \cdot a^b = a^{2b}$$

the identity case  $\mathbf{u}^H = \mathbf{I}$  has the property that

[illegible][illegible]

$$f(x) = \begin{pmatrix} x^2 & x^3 & x^4 \end{pmatrix} \quad \text{and} \quad g(x) = \begin{pmatrix} x^2 & x^3 & x^4 \end{pmatrix}$$

27476 CONCENTRATION OF THE CONCENTRATED PROPORTION TO PROPORTION 1 SPACE

[illegible]

$$\text{SE}(x^*) = I(x^*)^{-1} = \frac{1}{n} - \text{SE}(x^*)^2 = 0$$

REF ID: A66779

IN THE CASE OF THE "INDEPENDENT" REPORT, "25 JUNE."

SECRET//UNCLASSIFIED//FOR EYES OF THE PRESIDENT AND SELECT COMMITTEES

ON 21-20 SUBJ: CATH. BUS. ASING OF 21-20 OF CONCOMITANT FORCE.

TYPE: GRADE 2; SUBCATEGORIES: GRADATION OF 5 YEARS; TROUBLESHOOTING

on some of the things to be done on the 2<sup>nd</sup> of September.

It is clear from Definition 1 that its main domain is that of exponential family distributions, i.e. in the i.i.d. case distributions with densities of the form

$$f_{Y_j}(y|\theta) = \exp\{-\theta'b(y) + c(\theta) + d(y)\} \quad (3.7)$$

where  $\theta$ , possibly a vector, is unrestricted. Some further comment will be given on this later. It is not hard to see that consistency properties of the type (3.5) and (3.6) will hold generally for such exponential families. Indeed, the simple relationship between the sufficient statistic  $s = \sum b(y_j)$  and the maximum likelihood estimate  $\hat{\theta}$  for (3.7) enables one to show that, in obvious notation,

$$\text{PLIK}(s_F|s_P) = \exp\{-\hat{\theta}'_P s_F + n_F c(\hat{\theta}_P) + d^*(s_F)\} + o(n_P^{-1}). \quad (3.8)$$

Some details are given in the Appendix.

Before continuing with the general development of Definition 1, it is of interest to consider the inferential use of PLIK. According to the discussion of Section 2.2, this new likelihood would be used to remove non-uniqueness of confidence regions by choosing a region  $P_\alpha(y)$  such that

$$\sup_{z \in P_\alpha(y)} \text{PLIK}(s_F(z)|s_P(y)) \leq \inf_{z \in P_\alpha(y)} \text{PLIK}(s_F(z)|s_P(y)). \quad (3.9)$$

This would seem to be a non-controversial proposal. Now recall the other difficulty with critical region construction of confidence sets, namely the possibility that the hypothesis testing problem does not possess Neyman structure. In certain cases, typified by the following example, a controversial application of PLIK can overcome this difficulty, essentially by treating PLIK as a density.

It is clear from Definition 1 that the main theorem is that of exponential family distributions, 2.2. In the 2.2.5. case distributions which densities of the form

$$(2.7) \quad f_{\theta}(y) = \exp \{ \eta(\theta)'T(y) + \psi(\theta) \} h(y)$$

where  $\theta$  is possibly a vector, is understood. Some further comments will be given on this later. It is not hard to see that constancy properties of the type (2.7) and (2.8) will hold generally for such exponential families. Indeed, the relationship between the likelihood statistic  $s = \eta(\theta)$  and the maximum likelihood estimate for (2.7) enables one to show that, in obvious notation,

$$(2.8) \quad \eta(\hat{\theta}) = \exp \{ \eta(\theta)'T(\hat{y}) + \psi(\theta) \} h(\hat{y}) = \exp \{ \eta(\theta)'T(\hat{y}) + \psi(\theta) \} h(\hat{y})$$

Some details are given in the Appendix.

Before continuing with the formal development of Definition 1, it is of interest to consider the inferential use of ERM. According to the discussion of Section 2.2, this new likelihood would be used to remove non-uniqueness of confidence regions by choosing a region  $R(\gamma)$  such that

$$(2.9) \quad \inf_{\theta \in R(\gamma)} E_{\theta} [ \eta(\theta)'T(\hat{y}) + \psi(\theta) ] \geq \inf_{\theta \in R(\gamma)} E_{\theta} [ \eta(\theta)'T(\hat{y}) + \psi(\theta) ]$$

This would seem to be a non-conservative proposal. However, the other difficulty with official region construction of confidence sets, namely the possibility that the hypothesis testing problem does not possess a convex structure. In certain cases, verified by the following example, a conservative application of ERM can overcome this difficulty, essentially by treating ERM as a demand.

Example 3.2 Suppose that  $Y_1, \dots, Y_n$  are counts in non-overlapping unit time intervals of a Poisson process with constant rate  $\theta$ , and let  $Y_t^*$  be the corresponding count in a future interval of length  $t$ . Application of (3.3) easily shows that with  $S = \sum Y_j$

$$\text{PLIK}(y^*|s) = \frac{\Gamma(s + y^* + 1)}{\Gamma(s + 1)} n^s (n + t)^{s+y^*}. \quad (3.10)$$

But now suppose that we wish to make inference about the time interval  $T$  from the present to the occurrence of the next event; note that neither of the sampling-theory approaches of Section 2.2 could deal with this problem exactly. We observe that, in the earlier notation,  $T \geq t$  if and only if  $Y_t^* = 0$ . Therefore, the predictive likelihood of the event " $T \geq t$ " is precisely (3.10) evaluated at  $y^* = 0$ . It is then tempting to take

$$\text{PLIK}(t|s) = - \frac{\partial}{\partial t} \text{PLIK}(0|s) = sn^s (n + t)^{-s-1}, \quad (3.11)$$

although there is no logical case for this manipulation. Notice that the result satisfies the natural consistency property

$$\lim_{n \rightarrow \infty} \text{PLIK}(t|S) = \theta \exp(-\theta t) \quad \text{w.p.1.}$$

Not surprisingly, (3.11) is slightly different to the analogous (logical) result for inter-event time observations; in particular, for total time  $n$  between the start of observation and occurrence of the  $s^{\text{th}}$  event, with  $s$  fixed,

$$\text{PLIK}(t|n) = sn^{s-1} (n + t)^{-s}. \quad \square \square \square$$

Similar manipulations are possible in several other problems where the observable  $Y$  and future variable  $Z$  are determined by different sampling

ОБЪЕДИНЕНИЕ ИЛИ ИНОЕ ПРАВОПРИЧАСТНОЕ К ОБОИМ СЛУЖАВШИМ НА СЛУЖБЕ РАБОТНИКОВ  
ОБЪЕДИНЕНИЯ РАБОТНИКОВ ИЛИ ИНОЕ ПРАВОПРИЧАСТНОЕ К ОБОИМ СЛУЖАВШИМ НА СЛУЖБЕ РАБОТНИКОВ

$$\text{EPI}(C|U) = \sum_{u \in U} \frac{1}{|U|} \log \frac{1}{|C|}$$

APR 2 1964

КОД КОДЕТ ПУТЕ И РАДАНЕИ ДУЕ АСАИ ОУ ОРЕДИЛАНУОИ ДУЕ ОРЕДИКАТЕ ОУ  
 (ПРЕДНОУ) КАСИТЕ КОД РАДИО-КАСИТЕ ПУТЕ ОРЕДИКАТЕ: ДУЕ АСАИ ОУ  
 КОД АСАИ ОУ (ПРЕДНОУ) КАСИТЕ КОД РАДИО-КАСИТЕ ПУТЕ ОРЕДИКАТЕ: ДУЕ АСАИ ОУ

$$f_{\text{eff}}^{\text{BIBD}}(G, \phi) = \text{BIBD}(-\phi) \quad \text{A.B.I.}$$

[illegible]

$$\text{ВРЖ}(s|s) = -\frac{2}{\pi} \quad \text{ВРЖ}(0|s) = \frac{2}{\pi} (s - s)_{-s-1}^{-s-1} \quad (3.17)$$

penetration of the

and  $\bar{A}$  is the characteristic function of  $A$ . It is therefore  
 easy to see that  $\bar{A} = 0$  if and only if  $A$  is the empty set.  
 Moreover,  $\bar{A} = 1$  if and only if  $A$  is the universal set.  
 It is also easy to see that  $\bar{A} = 1 - A$ .  
 The following theorem is a direct consequence of the definition of  
 the characteristic function.

$$\text{EPR}(Q_n | S) = \frac{1}{1 + (2^n - 1)} \cdot \frac{1}{2^n} \cdot (1 + 2^n) \cdot \text{EPR}(Q_n | S) \quad (3.10)$$

установлено, что (3.3) имеет место для  $\lambda = 1$   
 и по две соответствующих точки  $\mu$  в каждой точке  $\lambda$  отрезка  $\lambda$ .  
 Пусть  $\mu_1$  и  $\mu_2$  — две точки  $\mu$  в точке  $\lambda$  соответствующие  $\lambda$  и  $\mu$ .  
 Тогда  $\mu_1$  и  $\mu_2$  — две точки  $\mu$  в точке  $\lambda$  соответствующие  $\lambda$  и  $\mu$ .

rules. The full consequences of such a wide interpretation of predictive likelihood have not been investigated.

To return to general developments, we now consider an extension of Definition 1 that is applicable in the non-i.i.d. exponential family case. It is not hard to see that in order for (3.3) to give the desired result for  $Z$ , independent of  $\theta$ , we require only that

(i) the conditional distribution of  $Z$  given  $s_F$  and  $s_P$  be known independent of  $\theta$

(ii)  $s_{PF} = s^*(s_P, s_F)$  for some function  $s^*(\cdot)$

and (iii)  $f_{s_{PF}}(s_{PF}; \theta) = k(s_P, s_F) f_{s_P}(s_P; \theta) f_{s_F|s_P}(s_F|s_P; \theta)$  . (3.12)

Thus, in general, it is not necessary for  $s_F$  to be a sufficient reduction of  $Z$ , nor is it sufficient for  $s_P$  to be a sufficient reduction of  $Y$ ; see Examples 3.3 and 3.4 below. We therefore make the revised

Definition 2 If  $s_P$  and  $s_F$  are minimal reductions of  $Y$  and  $Z$  such that

(i) - (iii) in (3.12) are satisfied, with  $s_{PF}$  the minimal sufficient reduction of  $(Y, Z)$ , then the predictive likelihood of  $s_F$  is

$$\text{PLIK}(s_F|s_P) = f_{s_P|s_{PF}}(s_P(y)|s_{PF}(y,z)) \quad (3.13a)$$

and the predictive likelihood of  $z$  is

$$\text{PLIK}_Z(z|s_P) = f_{Z|s_F, s_P}(z|s_F, s_P) \text{PLIK}(s_F|s_P) . \quad (3.13b)$$

(The statistic  $s_F$ , which depends on  $s_P$ , is called a minimal necessary statistic.) Two technical points to note are, firstly, that the condition (ii) can be satisfied if  $s_P$  is a minimal totally sufficient statistic as defined by Lauritzen [14]; and, secondly, Definition 2 extends the domain of non-trivial predictive likelihood beyond exponential

свойства для системы из нелинейных функциональных соотношений  
 заданных на области  $\Omega$  (1.1): для некоторой функции  $u$   
 существует (1.1) сдвигается на  $u$  по минимальной функции  
 свойства. Для некоторой функции  $u$  по сдвигу  $u$  на  $u$   
 (для некоторой  $u$  и для сдвига  $u$  на  $u$  по минимальной функции).

$$\text{или } \text{или } (u|v) = \int_{\Omega} u^{\alpha} v^{\beta} (u, v) \text{ или } (u^{\alpha}|v^{\beta}) \quad (3.13a)$$

или для функциональных соотношений

$$\text{или } (u^{\alpha}|v^{\beta}) = \int_{\Omega} u^{\alpha} v^{\beta} (u, v) \quad (3.13b)$$

свойства для (1.1) или для функциональных соотношений

(1) - (1.1) по (3.13) для некоторой функции  $u$  по минимальной функции  
 соотношения. То  $u$  и  $u^{\alpha}$  являются минимальными функциями  $u$  и  $u^{\alpha}$  для  
 $u$ : для некоторых  $u$  и  $u^{\alpha}$  по (3.13) для некоторых  $u$  и  $u^{\alpha}$   
 соотношения  $u$  по (3.13) для некоторых  $u$  и  $u^{\alpha}$  по (3.13) для некоторых  
 $u$  и  $u^{\alpha}$  по (3.13) для некоторых  $u$  и  $u^{\alpha}$  по (3.13) для некоторых

$$\text{или } (1.1) \text{ или } (u^{\alpha}|v^{\beta}) = \int_{\Omega} u^{\alpha} v^{\beta} (u, v) \quad (3.13c)$$

$$(1.1) \text{ или } (u^{\alpha}|v^{\beta}) = \int_{\Omega} u^{\alpha} v^{\beta} (u, v)$$

или соотношения

$$(1.1) \text{ или } (u^{\alpha}|v^{\beta}) = \int_{\Omega} u^{\alpha} v^{\beta} (u, v)$$

или для  $u$  соотношения  $u$  по (3.13) для некоторых  $u$  и  $u^{\alpha}$

или. То  $u$  по (3.13) для некоторых  $u$  и  $u^{\alpha}$  по (3.13) для некоторых  
 соотношения  $u$  по (3.13) для некоторых  $u$  и  $u^{\alpha}$  по (3.13) для некоторых

или соотношения  $u$  по (3.13) для некоторых  $u$  и  $u^{\alpha}$  по (3.13) для некоторых  
 соотношения  $u$  по (3.13) для некоторых  $u$  и  $u^{\alpha}$  по (3.13) для некоторых

или. То  $u$  по (3.13) для некоторых  $u$  и  $u^{\alpha}$  по (3.13) для некоторых

families, albeit not much beyond. The following examples illustrate the above discussion.

Example 3.3 Let  $Y_1, \dots, Y_n$  and  $Z = Y_{n+1}$  be i.i.d. with uniform density on  $[0, \theta]$ . Here the minimal sufficient statistic based on  $Y_1, \dots, Y_n$  is the maximum  $Y_{(n,n)}$ , so that  $S_{PF} = \max(Z, Y_{(n,n)})$ . Clearly,  $S_P = Y_{(n,n)}$ , but not all of  $Z$  is required to define  $S_{PF}$ . In fact to satisfy the condition (ii) in (3.12) we may take

$$S_F = \begin{cases} 0 & Z \leq Y_{(n,n)} \\ Z & Z > Y_{(n,n)} \end{cases}.$$

A simple probability calculation shows that

$$\text{PLIK}(s_F | s_P) = \begin{cases} n(n+1)^{-1} & (s_F = 0) \\ n(n+1)^{-1} s_F^{-1} s_P^{n-1} & (s_F > 0) \end{cases},$$

and since

$$f_{Z|S_F, S_P}(z | s_P, s_F) = \begin{cases} s_P^{-1} & (s_F = 0) \\ \delta(z - s_F) & (s_F > 0) \end{cases}$$

we have by (3.13) that

$$\text{PLIK}_Z(z | s_P) = \begin{cases} n(n+1)^{-1} s_P^{-1} & (s_F = 0) \\ n(n+1)^{-1} s_F^{-1} s_P^{n-1} & (s_F > 0) \end{cases}.$$

This last result is, in fact, a continuous probability density; notice that the predictive likelihood attached to the event  $Z \leq Y_{(n,n)}$ , i.e.  $n(n+1)^{-1}$ , formally agrees with the frequency probability of this event. The results for prediction of  $n_F > 1$  variables are obvious generalizations. □ □ □





Example 3.4 Let  $\{Y_j\}$  be a stationary first-order binary Markov sequence, with

$$\text{PR}(Y_{j+1} = b | Y_j = a) = \theta_{ab} \quad (a, b = 0, 1) . \quad (3.13)$$

Suppose that  $y = (y_0, y_1, \dots, y_{n_P})$  is observed, and that  $Z = (Y_{n_P+1}, \dots, Y_{n_P+n_F})$  is to be predicted. For any connected sequence of variables the minimal sufficient statistic is the initial value plus the matrix of one-step transition frequencies. Thus, if we let  $\delta(u) = 1(u = 0), = 0(u \neq 0)$  and define the matrix  $\underline{M}(r, S)$  by

$$M_{ab}(r, S) = \sum_{j=r}^{S-1} \delta(Y_j - a) \delta(Y_{j+1} - b) \quad (a, b = 0, 1; S \geq r) ,$$

we may express the minimal sufficient reductions of  $Y$ ,  $Z$  and  $(Y, Z)$  as

$$(Y_0, \underline{M}(0, n_P)), (Y_{n_P+1}, \underline{M}(n_P+1, n_P+n_F)) \text{ and } (Y_0, \underline{M}(0, n_P+n_F))$$

respectively. Clearly the first two do not give the last, because

$$\underline{M}(0, n_P+n_F) = \underline{M}(0, n_P) + \underline{M}(n_P+1, n_P+n_F) + \begin{pmatrix} \delta(Y_{n_P})\delta(Y_{n_P+1}) & \delta(Y_{n_P})\delta(1-Y_{n_P+1}) \\ \delta(1-Y_{n_P})\delta(Y_{n_P+1}) & \delta(1-Y_{n_P})\delta(1-Y_{n_P+1}) \end{pmatrix}$$

Thus in order to satisfy condition (ii) in (3.11) we must take

$$s_P = (Y_0, \underline{M}(0, n_P), Y_{n_P}), \quad s_F = (Y_{n_P+1}, \underline{M}(n_P+1, n_P+n_F));$$

$s_P$  is minimal totally sufficient.

The ensuing calculation of predictive likelihood is algebraically complicated, although the basic results required are given by Whittle [20]. For the special case  $n_F = 1$  it is straightforward to show that (3.12) gives

$$\text{PLIK}(y_{n_P+1} = b | y_0, \underline{m}(0, n), y_n = a) = \frac{m_{ab}(0, n) + 1}{1 + \sum_{k=0}^n m_{ak}(0, n)} .$$

$$\text{BTK}(A^{B+1}) = p(A^{0,1}(0^1,1) \cdot A^{B+1}) = \frac{1}{A^{B+1}(0^1,1) - 1} \quad (3.12)$$

Здесь

для всех значений  $B = 1, 2, 3, \dots$  справедливы следующие соотношения (3.13)

связывающие между собой функции, определенные на множестве  $[0,1]$ :

Эти соотношения являются следствием того, что функции  $A^B$  не являются линейными.

$$A^B = (A^{0,1}(0^1,1) \cdot A^{B+1}) \quad A^B = (A^{B+1}(0^1,1) \cdot A^{B+1})$$

причем если в первом соотношении (3.13) заменить  $A^B$  на  $A^{B+1}$ , то получим

$$A^{B+1}(0^1,1) = A^{B+1}(0^1,1) \cdot A^{B+1}(0^1,1) \quad A^{B+1}(0^1,1) = A^{B+1}(0^1,1) \cdot A^{B+1}(0^1,1)$$

следовательно, функции  $A^B$  и  $A^{B+1}$  связаны следующими соотношениями:

$$A^{B+1}(0^1,1) = A^{B+1}(0^1,1) \cdot A^{B+1}(0^1,1) \quad A^{B+1}(0^1,1) = A^{B+1}(0^1,1) \cdot A^{B+1}(0^1,1)$$

то есть функции  $A^B$  и  $A^{B+1}$  являются функциями от  $A^B$  и  $A^{B+1}$  соответственно.

$$A^{B+1}(0^1,1) = \frac{1}{A^{B+1}(0^1,1) - 1} \quad A^{B+1}(0^1,1) = \frac{1}{A^{B+1}(0^1,1) - 1}$$

следовательно, функции  $A^B$  и  $A^{B+1}$  являются функциями от  $A^B$  и  $A^{B+1}$  соответственно.

Следовательно, функции  $A^B$  и  $A^{B+1}$  являются функциями от  $A^B$  и  $A^{B+1}$  соответственно.

Следовательно, функции  $A^B$  и  $A^{B+1}$  являются функциями от  $A^B$  и  $A^{B+1}$  соответственно.

Следовательно, функции  $A^B$  и  $A^{B+1}$  являются функциями от  $A^B$  и  $A^{B+1}$  соответственно.

Следовательно, функции  $A^B$  и  $A^{B+1}$  являются функциями от  $A^B$  и  $A^{B+1}$  соответственно.

$$A^{B+1}(0^1,1) = p(A^B) = q(A^B) \quad (A^B = 0^1,1) \quad (3.14)$$

Здесь

функции  $p$  и  $q$  являются функциями от  $A^B$  и  $A^{B+1}$  соответственно.

This corresponds exactly to (3.4); the data in row 1-a of  $\underline{m}(0,n)$  is ignored when  $y_n = a$ , no matter how much information exists to support  $\theta_{0b} = \theta_{1b}$ .  $\square\square\square$

Although the predictive likelihood of Definition 2 can clearly play a role in sampling-theory confidence set construction, and possibly provide a credibility measure as in Example 3.2, there remains the question of the relation with Bayesian prediction. Recall from (2.5) that if  $\theta$  is the realized value of a random variable  $\Theta$  with prior density  $p(\theta)$ , then

$$f_{S_F|S_P}(s_F|s_P) = \frac{\int f_{S_P|\Theta}(s_P|\theta) f_{S_F|S_P,\Theta}(s_F|s_P,\theta) p(\theta) d\theta}{\int f_{S_P|\Theta}(s_P|\theta) p(\theta) d\theta}.$$

Now if condition (iii) in (3.12) is satisfied, we have immediately that

$$f_{S_F|S_P}(s_F|s_P) = \frac{\{k(s_P, s_F)\}^{-1} \int f_{S_{PF}|\Theta}(s_{PF}|\theta) p(\theta) d\theta}{\int f_{S_P|\Theta}(s_P|\theta) p(\theta) d\theta} \propto \text{PLIK}(s_F|s_P) f_{S_{PF}}^0(s^*(s_P, s_F)), \quad (3.14a)$$

where  $f_S^0$  is the marginal prior density of  $S$ . In the case of independent exponential family variables, when  $s_{PF} = s_P + s_F$  in the natural scale, (3.13a) may be expressed as

$$f_{S_F|S_P}(s_F|s_P) \propto \text{PLIK}(s_F|s_P) f_{S_F}^0(s_F). \quad (3.14b)$$

These results are precisely the desired analogs of the parametric result (2.6). By the consistency property of PLIK, and a corresponding consistency property of  $s_F$  as  $n_F \rightarrow \infty$ , the parametric result is in fact the limiting case of (3.14).

It is, then, possible to obtain posterior predictive distributions by assigning prior probabilities directly to the random variables, rather than to parameters, and using a predictive likelihood; this

bounded region  $\Omega$  is non-empty and the boundary  $\partial\Omega$  is a  
 regular surface of dimension  $n-1$  and the boundary  $\partial\Omega$  is  
 a compact set. The boundary  $\partial\Omega$  is a compact set.

It is then possible to obtain a bounded region  $\Omega$  of the  
 case of (3.13).

Consider  $\Omega$  of  $\mathbb{R}^n$  as  $\mathbb{R}^n$  and the boundary  $\partial\Omega$  is a compact set  
 (3.14). The boundary  $\partial\Omega$  is a compact set and a compact set  
 which is a bounded region of  $\mathbb{R}^n$  and the boundary  $\partial\Omega$  is a compact set.

$$e^{\frac{1}{2}A} e^{\frac{1}{2}B} (e^{\frac{1}{2}A} e^{\frac{1}{2}B}) = \text{tr}(e^{\frac{1}{2}A} e^{\frac{1}{2}B}) e^{\frac{1}{2}A} (e^{\frac{1}{2}B}). \quad (3.15)$$

(3.15) may be expressed as

which is a compact set. When  $e^{\frac{1}{2}A} = e^{\frac{1}{2}B} = e^{\frac{1}{2}C}$  in the boundary  $\partial\Omega$   
 where  $e^{\frac{1}{2}C}$  is the boundary  $\partial\Omega$  of  $\mathbb{R}^n$ . In the case of (3.15)

$$e^{\frac{1}{2}A} e^{\frac{1}{2}B} (e^{\frac{1}{2}A} e^{\frac{1}{2}B}) = \frac{e^{\frac{1}{2}A} (e^{\frac{1}{2}B}) e^{\frac{1}{2}A}}{\text{tr}(e^{\frac{1}{2}A} e^{\frac{1}{2}B})} = \frac{e^{\frac{1}{2}A} (e^{\frac{1}{2}B}) e^{\frac{1}{2}A}}{\text{tr}(e^{\frac{1}{2}A} e^{\frac{1}{2}B})} \quad (3.16)$$

Now the boundary  $\partial\Omega$  is a compact set. The boundary  $\partial\Omega$  is a compact set.

$$e^{\frac{1}{2}A} e^{\frac{1}{2}B} (e^{\frac{1}{2}A} e^{\frac{1}{2}B}) = \frac{e^{\frac{1}{2}A} (e^{\frac{1}{2}B}) e^{\frac{1}{2}A}}{\text{tr}(e^{\frac{1}{2}A} e^{\frac{1}{2}B})} = \frac{e^{\frac{1}{2}A} (e^{\frac{1}{2}B}) e^{\frac{1}{2}A}}{\text{tr}(e^{\frac{1}{2}A} e^{\frac{1}{2}B})}.$$

then

is the boundary  $\partial\Omega$  of a compact set. The boundary  $\partial\Omega$  is a compact set  
 of the boundary  $\partial\Omega$  of a compact set. The boundary  $\partial\Omega$  is a compact set  
 which is a compact set. The boundary  $\partial\Omega$  is a compact set. The boundary  $\partial\Omega$  is a compact set  
 which is a compact set. The boundary  $\partial\Omega$  is a compact set. The boundary  $\partial\Omega$  is a compact set.

which is a compact set. The boundary  $\partial\Omega$  is a compact set.

When  $A = B$  in the boundary  $\partial\Omega$  of a compact set. The boundary  $\partial\Omega$  is a compact set.

which is a compact set. The boundary  $\partial\Omega$  is a compact set.

provided Definition 2 is non-trivial for the particular model. In a sense this is not new, since for exponential families the prior density  $P(\theta)$  may be obtained from the marginal prior density  $f_{S_{PF}}^O(s_{PF})$  by inverting a Laplace transform. Thus (3.14) is a re-expression in suggestive terms of a known result.

Curiosity prompts one to ask: for what prior distributions is the (normalized) predictive likelihood a posterior predictive density in familiar problems? By (3.14) it is clearly necessary that  $S_{PF}$  have a uniform marginal prior density. The corresponding prior densities on  $\Theta$  for some commonly-occurring models are given in Table 1. Of course, the results depend on the sampling rule, since the predictive likelihood does; in that sense predictive likelihood has slightly lower status than parametric likelihood. However, if you pose a silly question, you get a silly answer.

[Table 1 here]

### 3.4 Conditional Predictive Likelihood

The predictive likelihood of Definition 2 gives a non-trivial result only when sufficiency arguments produce reductions of  $Y$  and  $Z$ . Thus the definition is meaningful for unrestricted exponential families and not much further. This leaves a large number of models. However, in many important problems the minimal sufficient statistic,  $S$ , although of dimension  $n$ , can be expressed as  $(T, C)$  where  $C$  is an ancillary statistic and  $T$  often has the same dimension as  $\theta$ . Conventionally inference would then proceed conditionally with the value of  $C$  regarded as fixed; see [3, Section 2.3] and [1], for example. In such cases it is possible to obtain a meaningful predictive likelihood by conditioning.



Suppose that conditions (i) - (iii) of (3.12) hold with  $S_P$ ,  $S_F$  and  $S_{PF}$  minimal. Also suppose that  $S_\alpha = (T_\alpha, C_\alpha)$  with  $C_\alpha$  ancillary for  $\alpha = P, F, PF$ . Clearly,  $C_{PF}$  is a function of  $(C_P, C_F, T_P, T_F)$  involving an ancillary function of  $(T_P, T_F)$ , and  $T_{PF}$  is a function of  $(T_P, T_F)$ . The statistic  $C_P$  contains no information about  $\theta$ , and hence none about  $S_F$ , so we condition on  $C_P$  and make

Definition 3 Under the above conditions,

$$\begin{aligned} \text{PLIK}(s_F | s_P) &= f_{C_F}(c_F) f_{T_P | S_{PF}, C_P}(t_P | s^*(s_P, s_F), c_P) \\ &= f_{C_F}(c_F) \text{PLIK}_C(t_F | s_P, c_F) . \end{aligned} \quad (3.15)$$

When  $C_{PF}$  is null, this reduces to Definition 2.

The following familiar example illustrates this definition.

Example 3.5 Let  $Y_1, \dots, Y_N$  be i.i.d. with continuous p.d.f.  $g(y-\theta)$ , supported on the whole real line, with  $Y = (Y_1, \dots, Y_{n_P})$  and  $Z = (Y_{n_P+1}, \dots, Y_{n_P+n_F})$ ,  $N = n_P + n_F$ . Denote the ordered values of  $Y$ ,  $Z$ , and  $(Y, Z)$  by  $\{Y_{(\alpha j)}, j=1, \dots, n_\alpha\}$  with  $\alpha = P, F$ , and  $PF$  respectively. Then in general  $S_\alpha = \{Y_{(\alpha j)}\}$  and we may take

$$T_\alpha = n_\alpha^{-1} \sum Y_{(\alpha j)}, \quad C_\alpha = \{Y_{(\alpha, j+1)} - Y_{(\alpha, j)}, j=1, \dots, n_\alpha - 1\}$$

for  $\alpha = P, F, PF$ ; it is only necessary that  $T_\alpha$  be a location statistic.

Now it is well-known and easy to show that, for each  $\alpha$ ,

$$f_{T_\alpha | C_\alpha}(t_\alpha | c_\alpha; \theta) = \frac{\prod_{j=1}^{n_\alpha} g(y_{(\alpha j)} - \theta)}{\int \prod_{j=1}^{n_\alpha} g(y_{(\alpha j)} - t) dt}$$

and

$$f_{C_\alpha}(c_\alpha) = n_\alpha! \int \prod_{j=1}^{n_\alpha} g(y_{(\alpha j)} - t) dt .$$



$$C^{\alpha}(c) = \int_{\Gamma} \frac{1}{\lambda^{\alpha}} \frac{1}{\lambda^{\alpha} - c} \lambda^{\alpha} d\lambda$$

или

$$C^{\alpha}(c|c^{\beta}) = \frac{\int_{\Gamma} \frac{1}{\lambda^{\alpha}} \frac{1}{\lambda^{\alpha} - c} \lambda^{\alpha} d\lambda}{\int_{\Gamma} \frac{1}{\lambda^{\alpha}} \frac{1}{\lambda^{\alpha} - c^{\beta}} \lambda^{\alpha} d\lambda}$$

полагая в формуле (3.1)  $\lambda^{\alpha} = c^{\beta}$  и получая формулу (3.2) для  $C^{\alpha}(c|c^{\beta})$ .

Для  $\alpha = 1$  и  $\beta = 1$  формула (3.2) принимает вид

$$C^1(c|c^1) = \frac{1}{c^1} \frac{1}{c^1 - c} \quad \text{или} \quad C^1(c|c^1) = \frac{1}{c^1} \frac{1}{c^1 - c} \quad (3.3)$$

или  $C^1(c|c^1) = \frac{1}{c^1} \frac{1}{c^1 - c}$  или формула (3.3)

или  $C^1(c|c^1) = \frac{1}{c^1} \frac{1}{c^1 - c}$  или формула (3.3)

или  $C^1(c|c^1) = \frac{1}{c^1} \frac{1}{c^1 - c}$  или формула (3.3)

или  $C^1(c|c^1) = \frac{1}{c^1} \frac{1}{c^1 - c}$  или формула (3.3)

или  $C^1(c|c^1) = \frac{1}{c^1} \frac{1}{c^1 - c}$  или формула (3.3)

или  $C^1(c|c^1) = \frac{1}{c^1} \frac{1}{c^1 - c}$  или формула (3.3)

или  $C^1(c|c^1) = \frac{1}{c^1} \frac{1}{c^1 - c}$  или формула (3.3)

$$C^1(c|c^1) = \frac{1}{c^1} \frac{1}{c^1 - c}$$

$$C^1(c|c^1) = \frac{1}{c^1} \frac{1}{c^1 - c} \quad (3.4)$$

или  $C^1(c|c^1) = \frac{1}{c^1} \frac{1}{c^1 - c}$  или формула (3.3)

или  $C^1(c|c^1) = \frac{1}{c^1} \frac{1}{c^1 - c}$  или формула (3.3)

или  $C^1(c|c^1) = \frac{1}{c^1} \frac{1}{c^1 - c}$  или формула (3.3)

или  $C^1(c|c^1) = \frac{1}{c^1} \frac{1}{c^1 - c}$  или формула (3.3)

или  $C^1(c|c^1) = \frac{1}{c^1} \frac{1}{c^1 - c}$  или формула (3.3)

или  $C^1(c|c^1) = \frac{1}{c^1} \frac{1}{c^1 - c}$  или формула (3.3)

или  $C^1(c|c^1) = \frac{1}{c^1} \frac{1}{c^1 - c}$  или формула (3.3)

It then follows readily from Definition 3 that

$$\begin{aligned}
 \text{PLIK}(s_F | s_P) &= f_{C_F}(c_F) f_{T_P | C_P}(t_P | c_P) f_{T_F | C_F}(t_F | c_F) \div f_{T_{PF} | C_{PF}}(t_{PF} | c_{PF}) \\
 &= n_F! \frac{\int \left\{ \prod_{j=1}^N g(y_j - u) \right\} du}{\int \left\{ \prod_{k=1}^{n_P} g(y_k - v) \right\} dv} ; \quad (3.16)
 \end{aligned}$$

for the unordered vector  $Z$  the factor  $n_F!$  disappears.

Note that (3.16) is formally identical to the pivotal (hence fiducial) density. The importance of conditioning is easily understood in this case: essentially we are predicting the difference of averages  $T_P - T_F$ , which, given the configuration ancillaries  $C_P$  and  $C_F$ , has a known distribution with mean not necessarily equal to zero, for example. Our inference must surely use this distributional information.  $\square \square \square$

The full implications of the conditional predictive likelihood definition are not yet realized. However, two facts are reasonably clear, particularly in light of the last example. Firstly, unless  $C_{PF}$  is null (as in Section 3.3), the Bayes result (3.14) does not in general hold. Secondly, there is a very close relationship with pivotal inference (Section 2.2), in the sense that the conditional predictive likelihood seems to be derived solely from  $C_{PF}$ .

the following lemma to be formulated together with  $C^{\mathbb{H}}$ .

THEOREM (3.3). In the above case the corresponding bilinear form  $B^{\mathbb{H}}$  associated with  $C^{\mathbb{H}}$  is a real symmetric bilinear form on  $V^{\mathbb{H}}$  (see in section 3.2). The bilinear form  $B^{\mathbb{H}}$  is non-degenerate if and only if the bilinear form  $B^{\mathbb{H}}$  is non-degenerate. The bilinear form  $B^{\mathbb{H}}$  is non-degenerate if and only if the bilinear form  $B^{\mathbb{H}}$  is non-degenerate.

The bilinear form  $B^{\mathbb{H}}$  is non-degenerate if and only if the bilinear form  $B^{\mathbb{H}}$  is non-degenerate.  $\square \square \square$

THEOREM (3.4). In the above case the bilinear form  $B^{\mathbb{H}}$  is non-degenerate if and only if the bilinear form  $B^{\mathbb{H}}$  is non-degenerate.

$B^{\mathbb{H}} = B^{\mathbb{H}}$  where  $B^{\mathbb{H}}$  is the bilinear form associated with  $C^{\mathbb{H}}$  and  $C^{\mathbb{H}}$  is a bilinear form on  $V^{\mathbb{H}}$ . The bilinear form  $B^{\mathbb{H}}$  is non-degenerate if and only if the bilinear form  $B^{\mathbb{H}}$  is non-degenerate.

THEOREM (3.5). In the above case the bilinear form  $B^{\mathbb{H}}$  is non-degenerate if and only if the bilinear form  $B^{\mathbb{H}}$  is non-degenerate.

$$\begin{aligned} & \int_0^1 \lambda^k (1-\lambda) d\lambda \\ & = \frac{1}{k+1} \int_0^1 \lambda^k d\lambda \\ & = \frac{1}{k+1} \left[ \frac{\lambda^{k+1}}{k+1} \right]_0^1 \\ & = \frac{1}{(k+1)^2} \end{aligned} \quad (3.6)$$

$$B^{\mathbb{H}}(C^{\mathbb{H}}|C^{\mathbb{H}}) = \frac{1}{k+1} \int_0^1 \lambda^k (1-\lambda) d\lambda = \frac{1}{(k+1)^2}$$

THEOREM (3.6). In the above case the bilinear form  $B^{\mathbb{H}}$  is non-degenerate if and only if the bilinear form  $B^{\mathbb{H}}$  is non-degenerate.

#### 4. An Empirical Analog of Predictive Likelihood

Thus far the discussion has dealt entirely with prediction in a parametric framework, the connection between data and predictand being the unknown constant  $\theta$ . This classical approach yields smooth and precise inferential statements by concentrating sample information on the few defining parameters of a known distribution family. A possible disadvantage of the approach is lack of flexibility in model fitting or, more accurately, failure to mirror the pragmatic approaches of applied statistics. Some concern on this point has recently stimulated work on the cross-validatory approach to prediction and model-fitting, notably by Stone [18] and Geisser [12]. The principal idea is to compare a (point) predictor estimated from a data subset with the predictands in the remaining data. Much the same idea underlies our definition of predictive likelihood in Section 3.2, which suggests the possibility of deriving an empirical analog of predictive likelihood by data sub-sampling.

We consider here only the simple case of one-dimensional exchangeable variables and prediction of a single outcome. Suppose that  $Y_1, \dots, Y_{n+1}$  are exchangeable, that we observe  $Y_1 = y_1, \dots, Y_n = y_n$  and wish to predict  $Y_{n+1}$ . Now let the statistic  $s(y_1, \dots, y_n)$  be chosen as a single data summary, presumably subjectively, so that differences in  $s$  will reflect important changes in the data. In the parametric framework  $s(\cdot)$  would be dictated by the minimal sufficiency requirement, but that is not so here; indeed one might well choose two or more statistics  $s_1, s_2, \dots$ . Because of the assumed exchangeability of  $Y_1, \dots, Y_n$ ,  $s(\cdot)$  is taken to be permutation-invariant. For this situation it is easy to develop an empirical analog of (3.3), as follows.

[illegible][illegible]

THE BRITISH COUNCIL OF EDUCATION

Let  $D_j = \{y_k : k=1, \dots, n, k \neq j\}$ ,  $j=1, \dots, n$  and  $D_{jk} = \{y_i : i=1, \dots, n, i \neq j, i \neq k\}$ ,  $k \neq j$ . Then by assumption all  $n$  subsamples  $D_j$  are equally likely, and all  $(n-1)$  subsamples  $D_{jk}$  within  $D_j$  are equally likely. We now have, in obvious notation,  $n(n-1)$  pairs  $(s(D_{jk}), s(D_j))$  which correspond to  $(S_P, S_{PF})$  as defined in Section 3.3. Suppose that there are  $r_n$  distinct values  $t_1, \dots, t_{r_n}$  of  $s(D_{jk})$ ,  $c_n$  distinct values  $u_1, \dots, u_{c_n}$  of  $s(D_j)$ , and define the  $r_n \times c_n$  matrix  $\underline{m}$  to have  $(a, b)$  element equal to the frequency of  $(t_a, u_b)$  in the totality of  $\{(s(D_{jk}), s(D_j)), 1 \leq k \neq j \leq n\}$  of nested subsamples. Then the matrix  $\underline{m}$  gives an empirical analog of (3.3), namely the subsampling frequency

$$\text{FREQ}(S_P = t_a | S_{PF} = u_b; n_P = n-2, n_F = 1) = \frac{m_{ab}}{m_{+b}}, \quad (3.17)$$

where  $m_{+b} = \sum_{a=1}^{r_n} m_{ab}$ . Now since  $u_b$  and  $t_a$  determine the  $(n-1)$  st. value of  $y_j$ , say  $y_{ab}^* = y^*(t_a, u_b)$ , the right hand side of (3.17) may be denoted  $\tilde{m}(t_a, y_{ab}^*; n-2)$ . By extrapolation we have

**Definition 4** If (3.17) is equal to  $\tilde{m}(t_a, y_{ab}^*; n-2)$ , then the empirical predictive likelihood of the event  $Y_{n+1} = y^*$ , given that  $s(y_1, \dots, y_n) = s$ , is

$$\text{EPLIK}(y^* | s) = \tilde{m}(s, y^*; n), \quad (3.18)$$

with support equal to  $\{y_{ab}^* : a=1, \dots, r_n; b=1, \dots, c_n\}$ .

**Example 4.1** Let the data consist of exchangeable Bernoulli trials, each  $y_j$  being 0 or 1, and suppose that  $s(y_1, \dots, y_n) = \sum y_j = s$ . A small calculation shows that  $r_n = 3$ ,  $c_n = 2$ ,  $(t_1, t_2, t_3) = (s-2, s-1, s)$ ,  $(u_1, u_2) = (s-1, s)$  and

$$\underline{m} = \begin{pmatrix} s(s-1) & 0 \\ s(n-s) & s(n-s) \\ 0 & (n-s)(n-s-1) \end{pmatrix}.$$



Then (3.17) and (3.18) give

$$\text{EPLIK}(y^*|s) = \left(\frac{n-s+1}{n+1}\right)^{1-y^*} \left(\frac{s+1}{n+1}\right)^{y^*}, \quad y^* = 0, 1,$$

which is exactly the parametric result (3.4) with  $n_F = 1$ . This is not very surprising, since our assumptions almost imply the binomial model.  $\square \square \square$

Sampling two deep is necessary because in the definition (3.3) it is  $S_{PF}$  that varies. For general  $n_F \geq 1$  the subsampling is of subsets size  $n_P - 2n_F$  from subsets size  $n_P - n_F$ . The assumption of a particular form for  $s(\cdot)$  is very close to a parametric assumption in Example 4.1, but not in general; for example  $\sum y_j$  is sufficient for several continuous families.

Example 4.2 Suppose that  $y_1, \dots, y_n$  are distinct values, and let  $s(y_1, \dots, y_n) = y_{(n)}$ , the largest order statistic, which is the lowest known upper bound on observations. Then  $r_n = 3$ ,  $c_n = 2$ ,  $(t_1, t_2, t_3) = (y_{(n-2)}, y_{(n-1)}, y_{(n)})$ ,  $(u_1, u_2) = (y_{(n-1)}, y_{(n)})$  and

$$\tilde{m} = \begin{pmatrix} 1 & 0 \\ n-2 & n-1 \\ 0 & (n-1)(n-2) \end{pmatrix}$$

For prediction of one further outcome  $y^*$ , (3.17) and (3.18) give

$$\text{EPLIK}(y^*|y_{(n)}) = \begin{cases} 1 - (n+1)^{-1}, & y^* \leq y_{(n)} \\ (n+1)^{-1}, & y^* > y_{(n)} \end{cases}.$$

This result is to be compared with the result for the uniform case in Example 3.3: here the likelihood beyond  $y_{(n)}$  is not distributed, since no density is assumed.  $\square \square \square$



Many interesting situations, such as that involving  $s = \sum y_j$ , encounter the difficulty that subsampled values of  $s$  may be distinct, leading to a matrix  $\underline{m}$  consisting of zeros and ones. If the empirical likelihood were to be generally useful, some smoothing device (acting on  $y$  or  $\underline{m}$ ) would be required. Our interest here is solely in noting the possibility of the cross-validation analog of Definition 1.

## 5. Concluding Remarks

In Section 2 some negative aspects of conventional sampling-theory prediction regions are emphasized, in order to give point to the subsequent discussion of likelihood. Clearly the sampling-theory methods are useful, and sensible confidence regions have usually been proposed that would agree with ordering by predictive likelihood PLIK. Some interesting work on prediction confidence regions may be found in [19], [6], [2], [13]; the latter is an exposition of invariance theory in pivotal prediction.

The treatment of predictive likelihood in Sections 3.3-4 clearly ignores some very interesting technical problems connected with the structure of sufficient statistics. There is currently much interest in such problems, particularly among Scandinavian statisticians; see particularly Lauritzen's work in [14], [15].

## 6. Acknowledgements

I particularly thank Steffen Lauritzen and Stephen Stigler for critical comments on an earlier draft of this article; Lauritzen suggested that the result (3.8) is true. Much of the work reported here was supported by NSF Grant No. NSF-MPS75-08778.



TABLE 1

Prior measures for which  $PLIK \propto$  Bayes posterior density

Model	Parameter $\theta$	Prior Measure
i.i.d. $MN(\mu, \Sigma)$ in $q$ dimensions same, $q=1, \Sigma=\sigma^2$	$\mu, \Sigma^{-1}$  $\mu, \sigma^2$	$ \Sigma ^{\frac{1}{2}q} d\mu d\Sigma^{-1}$  $\sigma d\mu d\sigma$
i.i.d. Bernoulli with $\text{pr}(Y=1)=\theta=\text{pr}(Y=0)$		
(a) fixed number of trials	$\theta$	$d\theta$
(b) fixed number of ones	$\theta$	$\theta^{-2} d\theta$
Poisson process with constant event rate $\lambda$		
(a) i.i.d. counts	$\lambda$	$d\lambda$
(b) i.i.d. inter-event times	$\lambda$	$\lambda^{-2} d\lambda$
i.i.d. Uniform on $[0, \theta]$	$\theta$	$\theta^{-1} d\theta$

## REFERENCES

- [1] Barnard, G.A. (1975). New foundations of statistical inference. Unpublished lecture notes, University of Minnesota.
- [2] Cox, D.R. (1974). Prediction intervals and empirical Bayes confidence intervals. (Unpublished).
- [3] Cox, D.R. and Hinkley, D.V. (1974). Theoretical Statistics. London: Chapman - Hall.
- [4] Edwards, A.W.F. (1974). A problem in the doctrine of chances. Proc. Conference on Foundational Questions in Statistical Inference, eds. Barndorff-Nielsen, O. et al. University of Aarhus.
- [5] Efron, B. (1975). Defining the curvature of a statistical problem (with applications to second order efficiency) Ann. Stat. 3, (to appear).
- [6] Faulkenberry, G.D. (1973). A method of obtaining prediction intervals. J. Amer. Statist. Assoc. 68, 433-435.
- [7] de Finetti, B. (1937). La Prévision: ses lois logiques, ses sources subjectives. Ann. Inst. H. Poincaré, 7, 1-68.
- [8] Fisher, R.A. (1922). On the mathematical foundations of theoretical statistics. Phil. Trans. Roy. Soc. London Ser. A 222, 309-368.
- [9] Fisher, R.A. (1935). The fiducial argument in statistical inference. Ann. Eug. 6, 391-398.
- [10] Fisher, R.A. (1973). Statistical Methods and Scientific Inference. (3rd Ed.) New York: Hafner.
- [11] Geisser, S. (1971). The inferential use of predictive distributions. In Foundations of Statistical Inference (V.P. Godambe and D.A. Sprott, eds.) pp. 456-469. Toronto, Montreal: Holt, Rinehart and Winston.
- [12] Geisser, S. (1975). The predictive sample reuse method with applications. J. Amer. Statist. Assoc. 79, (to appear).
- [13] Hora, R.B. and Buehler, R.J. (1967). Fiducial theory and invariant Prediction. Ann. Math. Statist. 38, 795-801.
- [14] Lauritzen, S.L. (1974). Sufficiency, prediction and extreme models. Scand. J. Statist. 1, 128-134.
- [15] Lauritzen, S.L. (1975). General exponential models for discrete observations. Scand. J. Statist. 2, 23-33.

[2] "Democracy" & P. (1980). "Democratization: A New Era of Global Change"

104. REFERENCES: A.P. (1940). "CONFIDENTIAL: INFORMATION AND ANALYSIS"

[REDACTED] (1981) [REDACTED] [REDACTED]

1781 received 2 (1912) the baggages and the other baggage in

[17] Congress: 8 (1911) The International Association of Agricultural

100-443887-1000

101. Report #7 (1639) - The document contains no information

[2] ITANON: NY (1933). ON THE NEUTRONIZATION REACTIONS OF

ALL TO HUNTER: R. (1221). To Squawpet: see above reference. see

[a] "Братская могила" с/п (ТД13) - в составе ее организационной структуры

[illegible][illegible][illegible]

(31) (Cox v. BIA) (1910) - Association: Supreme and ordinary police

[1] Journal 107 (1912): 108. Source: Journal of the American Medical Association.

- [16] Pearson, K. (1920). The fundamental problem of practical statistics. Biometrika 13, 1-16.
- [17] Pearson, K. (1921). Note on the 'fundamental problem of practical statistics'. Biometrika 16, 190-193.
- [18] Stone, M. (1974). Cross-validatory choice and assessment of statistical predictions, (with discussion). J. Roy. Statist. Soc. B 36, 111-147.
- [19] Weiss, L. (1955). A note on confidence sets for random variables. Ann. Math. Statist. 26, 142-144.
- [20] Whittle, P. (1955). Some distribution and moment formulae for the Markov chain. J. Roy. Statist. Soc. B 17, 235-242.

### Appendix: Consistency of PLIK

In Section 3.2 we refer without proof to two general consistency properties of PLIK in the i.i.d. exponential family case. A stronger result is stated without proof in (3.8). The lemma below and its corollaries give proofs of those results. For simplicity  $\theta$  is taken to be one-dimensional.

Lemma Let  $Y_1, \dots, Y_{n_{PF}}$  be i.i.d. with p.d.f.

$$f_Y(y; \theta) = \exp\{-\theta y + c(\theta) + d(y)\} ,$$

$S_P = \sum_{j=1}^{n_P} Y_j$ ,  $S_F = \sum_{j=1}^{n_F} Y_j$ ,  $S_{PF} = S_P + S_F$ ,  $n_{PF} = n_P + n_F$ . Then with  $n_F$  fixed,  $n_P \rightarrow \infty$ ,

$$f_{S_P|S_{PF}}(t|s+t) = f_{S_F|S_{PF}}(s|s+t) = f_{S_F}(t; \hat{\theta}_P) + o_p(n_P^{-1}), \quad (A.1)$$

where  $\hat{\theta}_\alpha$  is the maximum likelihood estimate determined by  $S_\alpha$  ( $\alpha = P, F, PF$ ).

Corollary Under the conditions of the Lemma, with PLIK given by

Definition 1,

$$(a) \quad \text{as } n_P \rightarrow \infty \quad \text{PLIK}(S_F|S_P) = f_{S_F}(s_F; \theta) + o_p(n_P^{-\frac{1}{2}})$$

$$(b) \quad \text{as } n_F \rightarrow \infty \quad \text{PLIK}(S_F|S_P) = f_{S_P}(s_P; \theta) + o_p(n_F^{-\frac{1}{2}})$$

Proof of Lemma The p.d.f. of  $S_\alpha$  may be denoted

$$g_{n_\alpha}(s_\alpha; \theta) = D(s_\alpha, n_\alpha) \exp\{-\theta s_\alpha + n_\alpha c(\theta)\} \quad (\alpha = P, F, PF).$$

Then (A.1) states that

$$f_{S_P|S_P+S_F}(s_P|s_P+t) = D(t, n_F) \exp\{-\hat{\theta}_P t + n_F c(\hat{\theta}_P)\} , \quad (A.2)$$



where  $S_\alpha + n_\alpha c'(\hat{\theta}_\alpha) = 0$  (A.3)

uniquely defines  $\hat{\theta}_\alpha$ . By definition

$$f_{S_P|S_P+S_F}(s|s+t) = f_{S_F|S_P+S_F}(t|s+t) = \frac{D(s, n_P)D(t, n_F)}{D(s+t, n_P+n_F)},$$

which can be re-expressed as

$$f_{S_P|S_P+S_F}(s|s+t) = D(t, n_F) \frac{g_{n_P}(s; \hat{\theta}_P)}{g_{n_{PF}}(s+t; \hat{\theta}_{PF})} \exp\{-n_P c(\hat{\theta}_P) + n_{PF} c(\hat{\theta}_{PF}) + \hat{\theta}_P s - \hat{\theta}_{PF}(s+t)\}. \quad (A.4)$$

But from (A.3) we deduce, by expansion, that

$$\hat{\theta}_P - \hat{\theta}_{PF} = \frac{n_P(s+t) - n_{PF}s}{n_P n_{PF} \{c''(\theta) + O_P(n_P^{-1/2})\}} = O_P(n_P^{-1}), \quad (A.5)$$

so that (A.4) becomes

$$f_{S_P|S_P+S_F}(s|s+t) = g_{n_F}(t; \hat{\theta}_P) \{1 + O_P(n_P^{-1})\} \frac{g_{n_P}(s; \hat{\theta}_P)}{g_{n_{PF}}\{s + O(1); \hat{\theta}_P + O_P(n_P^{-1})\}}.$$

It therefore remains to show that

$$\frac{g_{n_P}(s; \hat{\theta}_P)}{g_{n_{PF}}\{s + O(1); \hat{\theta}_P + O_P(n_P^{-1})\}} = 1 + O_P(n_P^{-1}).$$

This follows by applying (i)  $f_{X+O_P(n_P^{-1})}(x) = f_X(x) + O(n_P^{-1})$ ;

(ii)  $g_n(s + O(1); \theta) = g_n(s; \theta) + O(n^{-1})$ , (iii)  $g_n(s; \theta + O(n^{-1})) = g_n(s; \theta) + O(n^{-1})$ .

Proof of Corollary The first result (a) follows directly from (A.1) and

(iii) above because  $\hat{\theta}_P = \theta + O_P(n_P^{-1/2})$ . The second result (b) is a consequence of (A.1) with the roles of  $(S_P, n_P)$  and  $(S_F, n_F)$  interchanged.

The above arguments extend directly to the finite-dimensional parameter case, and to the exponential family linear model case. It is conjectured that corresponding results hold for Definitions 2 and 3, but this has not been proved.

Then (3.17) and (3.18) give

$$\text{EPLIK}(y^*|s) = \left(\frac{n-s+1}{n+1}\right)^{1-y^*} \left(\frac{s+1}{n+1}\right)^{y^*}, \quad y^* = 0, 1,$$

which is exactly the parametric result (3.4) with  $n_F = 1$ . This is not very surprising, since our assumptions almost imply the binomial model.  $\square \square \square$

Sampling two deep is necessary because in the definition (3.3) it is  $S_{PF}$  that varies. For general  $n_F \geq 1$  the subsampling is of subsets size  $n_P - 2n_F$  from subsets size  $n_P - n_F$ . The assumption of a particular form for  $s(\cdot)$  is very close to a parametric assumption in Example 4.1, but not in general; for example  $\sum y_j$  is sufficient for several continuous families.

Example 4.2 Suppose that  $y_1, \dots, y_n$  are distinct values, and let  $s(y_1, \dots, y_n) = y_{(n)}$ , the largest order statistic, which is the lowest known upper bound on observations. Then  $r_n = 3$ ,  $c_n = 2$ ,  $(t_1, t_2, t_3) = (y_{(n-2)}, y_{(n-1)}, y_{(n)})$ ,  $(u_1, u_2) = (y_{(n-1)}, y_{(n)})$  and

$$\tilde{m} = \begin{pmatrix} 1 & 0 \\ n-2 & n-1 \\ 0 & (n-1)(n-2) \end{pmatrix}$$

For prediction of one further outcome  $y^*$ , (3.17) and (3.18) give

$$\text{EPLIK}(y^*|y_{(n)}) = \begin{cases} 1 - (n+1)^{-1}, & y^* \leq y_{(n)} \\ (n+1)^{-1}, & y^* > y_{(n)} \end{cases}.$$

This result is to be compared with the result for the uniform case in Example 3.3: here the likelihood beyond  $y_{(n)}$  is not distributed, since no density is assumed.  $\square \square \square$

$$\tilde{H} = \begin{pmatrix} 0 & (N-1)(N-2) \\ N-2 & N-1 \\ 1 & 0 \end{pmatrix}$$

$$= (\lambda^{(N-2)} \cdot \lambda^{(N-1)} \cdot \lambda^{(N)}) \cdot (\lambda^1 \cdot \lambda^2) = (\lambda^{(N-1)} \cdot \lambda^{(N)}) \text{ and}$$

known from the theory of representations. When  $N = 3$ ,  $\lambda^1 = \lambda^2 = \lambda^3 = \lambda^4 = \lambda^5 = \lambda^6$

$\lambda^1 \cdot \lambda^2 \cdot \lambda^3 = \lambda^{(3)}$  and  $\lambda^4 \cdot \lambda^5 \cdot \lambda^6 = \lambda^{(3)}$  and the product

$\lambda^1 \cdot \lambda^2 \cdot \lambda^3 \cdot \lambda^4 \cdot \lambda^5 \cdot \lambda^6 = \lambda^{(6)}$  and the product

is

and the product is  $\lambda^{(6)}$  and the product is  $\lambda^{(6)}$

and the product is  $\lambda^{(6)}$  and the product is  $\lambda^{(6)}$

and the product is  $\lambda^{(6)}$  and the product is  $\lambda^{(6)}$

and the product is  $\lambda^{(6)}$  and the product is  $\lambda^{(6)}$

and the product is  $\lambda^{(6)}$  and the product is  $\lambda^{(6)}$

and the product is  $\lambda^{(6)}$  and the product is  $\lambda^{(6)}$

and the product is  $\lambda^{(6)}$  and the product is  $\lambda^{(6)}$

$$\text{rank}(A|e) = \left( \frac{N-1}{N-2} \right)_{1,1} \left( \frac{N-1}{N-2} \right)_{2,2} \cdot \lambda = 0 \cdot 1$$

and (3.10) and (3.11) are

Many interesting situations, such as that involving  $s = \sum y_j$ , encounter the difficulty that subsampled values of  $s$  may be distinct, leading to a matrix  $\underline{m}$  consisting of zeros and ones. If the empirical likelihood were to be generally useful, some smoothing device (acting on  $y$  or  $\underline{m}$ ) would be required. Our interest here is solely in noting the possibility of the cross-validation analog of Definition 1.